DETEKSI INFORMASI SENSITIF DALAM DOKUMEN TEKS DI SEKTOR JASA KEUANGAN DENGAN MODEL CNN BERBASIS TF-IDF DAN TF-RF

p-ISSN: 2355-7699

e-ISSN: 2528-6579

Steven Adriandi Vodegel*1, Septian Charles Vodegel2, Imelda3

1,2,3Fakultas Teknologi Informasi, Universitas Budi Luhur, Jakarta Email: ¹stevenvodegel97@gmail.com, ²septianvodegel@gmail.com, ³imelda@budiluhur.ac.id *Penulis Korespondensi

(Naskah masuk: 10 Oktober 2024, diterima untuk diterbitkan: 29 Oktober 2025)

Abstrak

Penelitian ini berfokus pada pengembangan model pembelajaran mesin untuk mendeteksi informasi sensitif dalam dokumen teks di industri jasa keuangan. Masalah utama yang diidentifikasi adalah potensi penyalahgunaan informasi oleh karyawan yang mengundurkan diri, keterbatasan metode deteksi tradisional, dan kebutuhan akan model pembelajaran mesin yang efektif. Ruang lingkup penelitian mencakup pengembangan model *Convolutional Neural Networks* (CNN) dengan metode pembebotan *Term Frequency-Inverse Document Frequency* (TF-IDF) dan *Term Frequency-Relevance Frequency* (TF-RF). Penelitian menggunakan pendekatan kuantitatif dan eksperimental, dengan tahapan pengumpulan data, pra-pemrosesan, penerapan pembebotan, pelatihan dan evaluasi model, serta validasi hasil. Data terdiri dari dokumen teks perusahaan jasa keuangan seperti laporan keuangan dan data nasabah. Pra-pemrosesan dilakukan untuk menghilangkan *noise* dan informasi tidak relevan, diikuti oleh metode pembebotan untuk memberi bebot pada kata-kata penting. Model CNN dilatih untuk mendeteksi pola yang menunjukkan informasi sensitif. Hasil penelitian menunjukkan metode TF-IDF lebih baik daripada TF-RF dalam mendeteksi informasi sensitif, dengan akurasi tertinggi 93,26%. Model CNN mampu mengenali pola kompleks dan mendeteksi informasi sensitif dengan akurasi tinggi. Evaluasi dengan akurasi, presisi, *recall*, dan *f1-score* menunjukkan bahwa model ini dapat diandalkan dan diaplikasikan dalam situasi nyata. Penelitian ini berkontribusi pada keamanan informasi dan penerapan pembebotan dalam meningkatkan kinerja model pembelajaran mesin.

Kata kunci: deteksi informasi sensitive, pembelajaran mesin, convolutional neural networks, term frequency-inverse document frequency, term frequency-relevance frequency.

SENSITIVE INFORMATION DETECTION IN TEXT DOCUMENTS WITHIN THE FINANCIAL SERVICES SECTOR USING A CNN MODEL BASED ON TF-IDF AND TF-RF

Abstract

This research focuses on the development of a machine learning model to detect sensitive information in text documents within the financial services industry. The main issues identified are the potential misuse of information by employees who resign, the limitations of traditional detection methods, and the need for an effective machine learning model. The scope of the research includes the development of a Convolutional Neural Networks (CNN) model with Term Frequency-Inverse Document Frequency (TF-IDF) and Term Frequency-Relevance Frequency (TF-RF) weighting methods. The study employs a quantitative and experimental approach, with stages including data collection, preprocessing, application of weighting methods, model training and evaluation, and result validation. The data consists of text documents from financial services companies such as financial reports and customer data. Preprocessing was carried out to remove noise and irrelevant information, followed by the application of weighting methods to assign importance to significant words. The CNN model was trained to detect patterns indicating sensitive information. The results show that the TF-IDF method performed better than TF-RF in detecting sensitive information, with the highest accuracy of 93.26%. The CNN model was able to recognize complex patterns and detect sensitive information with high accuracy. Evaluation using accuracy, precision, recall, and f1-score demonstrates that this model is reliable and applicable in real-world situations. This research contributes to information security and the use of weighting methods to improve the performance of machine learning models.

Keywords: sensitive information detection, machine learning, convolutional neural networks, term frequency-inverse document frequency, term frequency-relevance frequency.

1. PENDAHULUAN

Dalam era digital, informasi menjadi aset penting bagi perusahaan, khususnya di sektor jasa keuangan. Data sensitif seperti nasabah, strategi bisnis, dan laporan keuangan mendukung operasional perusahaan. Keamanan data ini krusial untuk menjaga reputasi, kepercayaan pelanggan, dan mematuhi regulasi. Namun, tantangan keamanan informasi terus meningkat seiring ancaman siber yang semakin kompleks, seperti yang dikemukakan oleh (Soldatos & Kyriazis, 2022), sehingga dibutuhkan solusi lebih canggih. Meskipun beragam metode dan alat telah dikembangkan, banyak yang kurang efektif karena implementasi yang sulit dan rendahnya pemahaman pengguna terhadap teknologi keamanan tersebut (Templ & Sariyar, 2022).

Ancaman keamanan ini tidak hanya berasal dari luar perusahaan, tetapi juga dari dalam. Salah satu contoh nyata adalah potensi penyalahgunaan informasi oleh karyawan yang akan mengundurkan diri (Al-shehari & Alsowail, 2021). Dalam satu kasus di perusahaan jasa keungan, seorang karyawan diduga mengunduh dokumen sensitif beberapa hari sebelum berhenti bekerja. Selain itu, ancaman lain seringkali disebabkan oleh kesalahan manusia, seperti tidak mengunci perangkat elektronik atau penggunaan flashdisk yang tidak aman (Ali, Jalal & 2020). Al-Obavdv. Ibrahem Peristiwa menunjukkan kelemahan dalam mendeteksi aktivitas mencurigakan, terutama yang berasal dari pihak internal, meskipun kebijakan keamanan sudah diterapkan. Hal ini mendukung temuan Shakti & Hidayanto (2024) bahwa kesadaran terhadap keamanan informasi masih rendah di kalangan pengguna.

Savangnya, metode tradisional dalam mendeteksi aktivitas mencurigakan sering kali tidak memadai. Volume data yang besar dan kompleksitas dalam mengidentifikasi informasi sensitif secara manual membuat pendekatan konvensional ini kurang efisien (Nugroho, Istiadi & Marisa, 2020). Oleh karena itu, dibutuhkan teknologi yang lebih mutakhir untuk mendeteksi informasi sensitif dan mencegah penyalahgunaannya. Salah satu solusi terbukti efektif adalah penggunaan pembelajaran mesin. Model pembelajaran mesin mampu mengenali pola-pola data yang rumit serta mendeteksi aktivitas mencurigakan dengan akurasi yang lebih tinggi. Penelitian oleh Ali et al. (2022) menunjukkan bahwa pembelajaran memberikan hasil signifikan dalam mendeteksi kecurangan keuangan dan mampu menangani volume data yang besar dengan efisiensi tinggi.

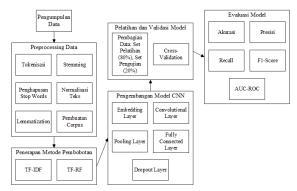
Dalam penelitian ini, fokus diarahkan pada perbandingan dua metode pembobotan, yaitu Term Frequency-Inverse Document Frequency (TF-IDF) dan Term Frequency-Relevance Frequency (TF-RF), yang sering digunakan dalam pengolahan data teks. Kedua metode ini memiliki keunggulan masingmasing dalam menangkap fitur relevan dari teks dan

memberikan kontribusi signifikan pada model pembelajaran mesin (Irena & Setiawan, 2020; Tama & Sibaroni, 2023). Dalam konteks industri jasa keuangan, pendekatan ini diterapkan untuk mendeteksi informasi sensitif yang berpotensi disalahgunakan.

Dalam konteks ini, Convolutional Neural Networks (CNN), salah satu jenis pembelajaran mesin yang sangat canggih dan popular, menawarkan kemampuan yang unggul dalam mendeteksi informasi sensitif dalam dokumen teks. CNN juga memberikan nilai positif untuk menjadi dasar model utama karena mampu menghasilkan performa yang baik melaui pembelajar fitur yang kuat (Mahardika, Yudistira & Ridok, 2024). CNN dalam struktur jaringan berlapis-lapisnya (layer-layer) memungkinkan proses pembelajaran berlangsung lebih efisien (Nurdiawan, Susana & Fagih, 2024). CNN telah terbukti efektif di berbagai aplikasi pemrosesan bahasa alami (Natural Language Processing), seperti klasifikasi teks, analisis sentimen, dan deteksi spam. Keuntungan utama CNN dibandingkan dengan pendahulunya kemampuannya untuk secara otomatis mendeteksi fitur-fitur penting tanpa pengawasan manusia, yang membuatnya menjadi yang paling banyak digunakan (Alzubaidi et al., 2021). CNN sebagai model dalam deep learning lebih baik digunakan pada penelitian ini karena kebutuhan jumlah parameter yang lebih sedikit sehingga komputasi lebih ringan serta keterbatasan dataset, tidak seperti RNN dan BERT yang membutuhkan parameter lebih besar serta harus menggunakan dataset lebih besar untuk menghindari overfitting.

Model ini dapat dilatih untuk mengenali pola spesifik yang mengindikasikan keberadaan informasi sensitif dalam dokumen dengan memanfaatkan metode pembobotan seperti Term Frequency-Inverse Document Frequency (TF-IDF) dan Term Frequency-Relevance Frequency (TF-RF). Penelitian yang dilakukan oleh Stojanović et al. (2021) juga mendukung penggunaan CNN dalam meningkatkan deteksi kecurangan di aplikasi financial technology (fintech). Penelitian oleh Yuhana et al. (2022) menunjukkan bahwa metode pembelajaran mendalam, khususnya CNN, RNN, dan HAN, efektif dalam klasifikasi teks untuk pembelajaran adaptif.

Pendekatan pembobotan TF-IDF memberikan keunggulan dalam menilai pentingnya suatu kata berdasarkan distribusinya di seluruh dokumen, sementara TF-RF menambahkan dimensi relevansi dengan mempertimbangkan frekuensi kemunculan di dokumen relevan (Irena & Setiawan, 2020; Tama & Sibaroni, 2023). Dengan memanfaatkan kedua metode ini, penelitian ini bertujuan untuk mengidentifikasi teknik terbaik yang dapat



Gambar 1. Tahapan Penelitan

mendukung performa CNN dalam deteksi informasi sensitif.

Penelitian yang dilalukan Suartana (2022) menunjukkan bahwa deep learning memiliki potensi besar dalam mendeteksi dan mengklasifikasikan serangan jaringan melalui analisis data skala besar dan pembelajaran pola informasi. Namun, penelitian ini membawa pendekatan deep learning lebih jauh dengan mengaplikasikan CNN yang didukung pembobotan TF-IDF dan TF-RF dalam konteks deteksi informasi sensitif pada dokumen teks, khususnya untuk mengatasi tantangan keamanan informasi di sektor jasa keuangan

Dengan mengembangkan dan mengevaluasi model pembelajaran mesin berbasis Convolutional Neural Networks (CNN) yang dipadukan dengan metode pembobotan Term Frequency-Inverse Document Frequency (TF-IDF) dan Frequency-Relevance Frequency (TF-RF), penelitian ini memberikan kontribusi baru dalam upaya otomatisasi deteksi informasi sensitif. Model ini dirancang untuk meningkatkan akurasi identifikasi potensi penyalahgunaan informasi oleh karyawan yang akan resign, memungkinkan perusahaan untuk lebih proaktif dalam mencegah kebocoran data. Solusi ini tidak hanya meningkatkan keamanan informasi perusahaan tetapi juga menawarkan pendekatan inovatif dalam pengolahan data yang seimbang dengan perlindungan privasi (Guha et al., 2021). Penelitian ini diharapkan dapat memberikan kontribusi signifikan dalam bidang keamanan informasi, khususnya dalam penerapan pembelajaran mesin untuk mendeteksi informasi sensitif secara Dengan demikian, otomatis. penelitian berkontribusi pada pengembangan metode keamanan informasi berbasis pembelajaran mesin, khususnya di sektor jasa keuangan.

2. METODE PENELITIAN

Gambar 1 menunjukkan tahapan penelitian meliputi proses inti dalam studi ini yang bertujuan untuk mencegah terjadinya penyalahgunaan informasi sensitif. Pada penelitian ini pembuatan model merujuk pada penelitian (Guha et al., 2021) vang menerapkan metode supervised learning dengan dua kelas/ label data atau biasa disebut binary

Tabel 1. Pengaturan Eksperimen & Informasi Arsitektur CNN

ID	Tipe Lapisan/ layer	Parameters
1	Optimzer	Adam
2	Loss Function	Binary Crossentropy
3	<i>Epochs</i>	10
4	Split dataset	80% Training, 20% Testing
5	Random State	42
6	Conv1D	Filters: 128, Kernel size: 5, Activation: ReLU
7	MaxPooling1D	Pool size: 2
8	Conv1D	Filters: 64, Kernel size: 5, Activation: ReLU
9	MaxPooling1D	Pool size: 2
10	Flatten	Yes
11	Dense	Units: 64, Activation: ReLU
12	Dropout	Rate: 0.5
13	Dense (output layer)	Units: 1, Activation: Sigmoid

Tabel 2. Dataset Penelitian

ID	Konten Kalimat	Label
1	Pemblokiran serta merta	sensitif
2	Berita keuangan	non-sensitif
3	Rincian infrastuktur TI	sensitif
4	Portofolio investasi	sensitif
5	Hedge fund	sensitif
6	Laporan industry	non-sensitif
7	Saldo akun	sensitif
8	Umpan balik nasabah	sensitif
9	Laporan kinerja divisi perencanaan strategis	non-sensitif
10	Nomor rekening	sensitif

classification. Proses dimulai dari pengumpulan data hingga implementasi model CNN untuk melakukan deteksi informasi sensitif. Pada penelitian ini untuk pengaturan eksperimen & informasi arsitektur CNN, dapat merujuk pada Tabel 1.

2.1. Pengumpulan Data

Obyek penelitian dalam tesis ini adalah dokumen teks yang berasal dari lingkungan perusahaan jasa keuangan seperti bank, asuransi, dan perusahaan investasi. Dataset yang digunakan mencakup berbagai jenis dokumen dengan label yang menunjukkan apakah dokumen tersebut mengandung informasi sensitif atau non-sensitif. Dataset ini dikumpulkan dari berbagai departemen dalam perusahaan jasa keuangan dan telah melalui proses pra-pemrosesan untuk memastikan kualitas dan relevansinya. Pelabelan dataset dilakukan secara manual dengan memasikan sesuai dengan peraturan internal perusahaan serta mempertimbangkan saran dari subject matter expert. Jenis dokumen meliputi data nasabah, laporan keuangan, strategi bisnis, dokumen internal perusahan dengan total 442 data. Pada Tabel 2. Menunjukkan dataset yang digunakan.

"Pemblokiran serta merta" (ID 1) dilabeli sebagai sensitif karena mungkin berisi informasi terkait tindakan keamanan mendadak yang harus dirahasiakan. Sebaliknya, "berita keuangan" (ID 2) dianggap non-sensitif karena informasi ini biasanya bersifat publik dan dapat diakses oleh semua orang. "Rincian infrastuktur TI" (ID 3) dikategorikan sensitif karena mencakup detail teknis yang dapat membahayakan infrastruktur perusahaan jika jatuh ke pihak yang salah. "Saldo akun" (ID 7) juga dianggap sensitif karena memuat informasi finansial pribadi yang sangat rentan terhadap penyalahgunaan. Sementara itu, "Laporan kinerja divisi perencanaan strategis" (ID 9) dilabeli sebagai non-sensitif karena laporan ini kemungkinan tidak mengandung informasi rahasia. Namun, jika laporan ini mencakup rencana masa depan yang strategis, bisa jadi diklasifikasikan sebagai sensitif.

Dalam pengembangan model membutuhkan langkah pra-pemprosesan data yang bersifat krusial (Ahsan et al., 2021). Langkah ini melibatkan beberapa tahapan penting, seperti tokenisasi, yang memecah teks menjadi unit-unit kecil seperti kata atau frasa untuk memudahkan analisis. Stemming dilakukan untuk mengubah kata ke bentuk dasarnya, yang membantu mengurangi variasi kata dan meningkatkan akurasi analisis. Selain itu. penghapusan stop words dilakukan untuk menghilangkan kata-kata umum yang memberikan informasi penting, sehingga analisis dapat fokus pada kata yang lebih signifikan. Normalisasi teks memastikan konsistensi dengan mengubah huruf kapital menjadi huruf kecil dan menghapus tanda baca yang tidak relevan. Lemmatization mengembalikan kata ke bentuk dasar dengan mempertimbangkan konteksnya. Seluruh langkah ini memastikan bahwa data yang digunakan dalam pelatihan model adalah data yang bersih, konsisten, dan siap dianalisis.

2.2. Penerapan Metode Pembobotan

Penerapan metode pembobotan merupakan langkah penting dalam pemrosesan data teks untuk memastikan bahwa model CNN dapat mengenali dan memprioritaskan informasi yang relevan. Dua metode pembobotan yang digunakan dalam penelitian ini adalah *Term Frequency-Inverse Document Frequency* (TF-IDF) dan *Term Frequency-Relevance Frequency* (TF-RF).

TF-IDF adalah metode pembobotan yang digunakan untuk mengevaluasi pentingnya sebuah kata dalam dokumen relatif terhadap *korpus* dokumen. *Term Frequency* (TF) mengukur seberapa sering sebuah kata muncul dalam dokumen. TF dari kata *t* dalam dokumen *d* dapat dirumuskan sebagai:

$$TF(t,d) = \frac{f_{t,d}}{N_d} \tag{1}$$

Dimana $f_{t,d}$ adalah frekuensi kemunculan kata t dalam dokumen d dan N_d adalah total jumlah kata dalam dokumen d. Sedangkan *Inverse Document Frequency* (IDF) mengukur seberapa sering sebuah kata muncul di seluruh *korpus* dokumen. IDF dari kata t dapat dirumuskan sebagai:

$$IDF(t) = \log\left(\frac{N}{DF(t)}\right)$$
 (2)

Dimana N adalah total jumlah dokumen dalam korpus, dan DF (t) adalah jumlah dokumen yang mengandung kata t. Kombinasi dari kedua komponen ini memberikan persamaan TF-IDF:

$$TF - IDR(t, d) = TF(t, d) x \log \left(\frac{N}{DF(t)}\right)$$
 (3)

TF-RF adalah metode pembobotan yang mempertimbangkan relevansi kata dalam konteks tertentu yang bertujuan memberikan performa pembobotan yang lebih baik disbanding metode lain (Harmandini & L, 2024). Metode ini dipilih juga karena dalam konteks klasifikasi teks, pembobotan TF RF mampu memberikan hasil yang optimal (Dananjaya & Indradewi, 2023; Sari et al., 2022; Ramadhan, Sari & Adikara, 2021). Metode ini memberikan bobot yang lebih tinggi pada kata-kata yang lebih relevan dengan konteks spesifik dokumen sensitif. Metode ini juga terdiri dari dua komponen utama: Term Frequency (TF) dan Relevance Frequency (RF). Relevance Frequency (RF) mengukur relevansi sebuah kata dalam dokumen dengan memperhatikan frekuensi kata dalam dokumen yang relevan. RF dari kata t dapat dirumuskan sebagai:

$$RF(t) = \frac{R}{r_t} \tag{4}$$

Dimana R adalah total jumlah dokumen relevan, dan r_t adalah jumlah dokumen relevan yang mengandung kata t. Kombinasi dari kedua komponen ini memberikan nilai TF-RF:

$$TF - RF(t,d) = TF(t,d) \times RF(t)$$
 (5)

2.3. Pelatihan dan Evaluasi Model CNN

Model CNN yang digunakan adalah model CNN yang dirancang manual menggunakan lapisanlapisan untuk memenuhi kebutuhan tertentu dan tidak menggunakan model CNN spesifik seperti VGG dan ResNet.. Library yang digunakan adalah keras. Model CNN dilatih menggunakan set pelatihan dan dievaluasi menggunakan set pengujian untuk mencapai performa terbaik. Dataset dibagi menjadi set pelatihan (80%) dan set pengujian (20%) menggunakan teknik k-fold cross-validation. Teknik ini membagi dataset menjadi beberapa subset, melatih model pada k-1 subset, dan menguji pada subset yang tersisa. Proses ini diulang k kali sehingga setiap data menjadi bagian dari set pengujian satu kali, memungkinkan model diuii secara menyeluruh serta menghindari bias dalam pembagian data. Hal ini sejalan dengan fakta bahwa k-fold cross-validation merupakan metode validasi yang andal, meskipun tidak digunakan sebanyak metode hold-out dalam bidang educational data mining (Ghorbani & Ghousi,

Model CNN yang dirancang memiliki beberapa lapisan utama, yaitu *embedding layer*, *convolutional*

layer, pooling layer, dan hidden layer (yang terdiri dari fully connected layer dan dropout layer) (Yuhana et al., 2022). Setiap lapisan memiliki peran spesifik dalam menganalisis dan mendeteksi pola dalam teks yang menunjukkan keberadaan informasi sensitif. Penelitian yang dilakukan oleh Rachman & Santoso (2021) menunjukkan bahwa metode-metode deep learning seperti CNN, RNN, LSTM, dan Bidirectional LSTM, bersama dengan teknik NLP seperti tokenization dan word embeddings, sangat efektif dalam klasifikasi teks untuk analisis sentimen.

elama pelatihan, model mengoptimalkan parameter melalui iterasi berulang hingga mencapai konvergensi, sementara teknik regularisasi seperti dropout dan batch normalization diterapkan untuk mencegah overfitting. Algoritma optimasi seperti Adam atau RMSprop dipakai untuk mempercepat konvergensi dan meningkatkan performa. Evaluasi model dilakukan menggunakan metrik seperti akurasi, presisi, recall, f1-score, dan AUC-ROC untuk menilai kemampuan model dalam mendeteksi informasi sensitif yang nilainya mengacu pada confusion matrix (Krisnabayu, Ridok & Budi, 2021). Adapun matrik yang digunakan mengikuti penelitian sebelumnya (Mahardika, Yudistira & Ridok, 2024) (Akbar et al., 2024), dimana nilai matrik ini dapat menunjukkan kemampuan model. Evaluasi ini memberikan gambaran kinerja model pada data yang belum pernah dilihat sebelumnya, memastikan performa yang baik secara keseluruhan.

$$akurasi = \frac{TP + TN}{TP + TN + FP + FN} \tag{6}$$

$$presisi = \frac{TP}{TP + FP} \tag{7}$$

$$f1 - score = \frac{2 x recall x precision}{recall + precision}$$
 (8)

Tajbakhsh et al. (2017) menyatakan pada penelitiannya bahwa model **CNN** dikembangkan dapat memerlukan penyesuaian atau fine-tuning untuk meningkatkan kinerjanya. Proses ini melibatkan penyesuaian parameter seperti learning rate, batch size, dan jumlah epoch menggunakan metode grid search atau random search untuk menemukan kombinasi terbaik. Selain itu, eksperimen dengan arsitektur model, seperti menambah lapisan atau mengubah ukuran kernel, dilakukan untuk menemukan konfigurasi optimal. Metode pembobotan juga dioptimalkan dengan meninjau kembali penggunaan TF-IDF dan TF-RF serta mencoba metode lain untuk meningkatkan akurasi dalam mendeteksi informasi sensitif.

2.4. Prosedur Keamanan dan Pengolahan Data

Untuk menjaga keamanan dan integritas data selama penelitian, beberapa prosedur keamanan diterapkan. Data yang disimpan di Google Drive dienkripsi menggunakan algoritma standar industri untuk melindunginya dari akses tidak sah, baik selama penyimpanan maupun transfer. Selain itu, akses ke data dan kode penelitian dibatasi hanya untuk individu yang berwenang, dengan autentikasi dua faktor (2-FA) diterapkan untuk meningkatkan keamanan. Setiap aktivitas akses juga diaudit secara berkala mendeteksi adanya guna Selain itu, pencadangan mencurigakan. dilakukan secara teratur di lokasi terpisah yang dilindungi dengan baik, dan dijadwalkan secara otomatis untuk memastikan salinan cadangan selalu up-to-date, sehingga menghindari risiko kehilangan data.

HASIL DAN PEMBAHASAN

3.1. Hasil Eksperimen Pertama TF-IDF

Pada eksperimen pertama TF-IDF, model yang digunakan untuk mengklasifikasikan data dievaluasi menggunakan beberapa metrik performa yaitu akurasi, presisi, recall, F1-score, dan AUC-ROC yang ditunjukkan pada Tabel 3. Akurasi model sebesar 0.876 menunjukkan bahwa model mampu memprediksi dengan benar sekitar 87.6% dari seluruh data uji yang digunakan. Presisi untuk kelas 0 adalah 0.86 dan untuk kelas 1 adalah 0.89, mengindikasikan bahwa dari seluruh prediksi positif yang dibuat oleh model, 86% untuk kelas 0 dan 89% untuk kelas 1 adalah benar. Recall untuk kelas 0 adalah 0.92 dan untuk kelas 1 adalah 0.83, yang menunjukkan bahwa model mampu menemukan 92% dari semua kasus positif sebenarnya untuk kelas 0 dan 83% untuk kelas 1. F1-score, yang merupakan rata-rata harmonis antara presisi dan recall, tercatat sebesar 0.89 untuk kelas 0 dan 0.86 untuk kelas 1. Nilai AUC-ROC model sebesar 0.901 menandakan kemampuan model yang sangat baik dalam membedakan antara kelas positif dan negatif. Semakin mendekati kurva ROC ke sudut kiri atas, semakin baik performa model dalam memisahkan kelas positif dari kelas negatif.

Pada eksperimen kedua TF-IDF, model yang digunakan untuk mengklasifikasikan data dievaluasi menggunakan metrik performa yang mencakup akurasi, presisi, recall, F1-score, dan AUC-ROC yang ditunjukkan pada Tabel 4. Akurasi model tercatat sebesar 0.932, yang menunjukkan bahwa model mampu memprediksi dengan benar sekitar 93.2% dari seluruh data uji yang digunakan. Presisi untuk kelas 0 dan kelas 1 sama-sama sebesar 0.93, yang mengindikasikan bahwa dari seluruh prediksi positif yang dibuat oleh model, 93% adalah benar untuk kedua kelas. Recall untuk kelas 0 dan kelas 1 juga sama-sama sebesar 0.93, yang menunjukkan bahwa model mampu menemukan 93% dari semua kasus positif yang sebenarnya untuk kedua kelas. F1score, yang merupakan rata-rata harmonis antara presisi dan recall, tercatat sebesar 0.93 untuk kedua kelas. Nilai AUC-ROC model sebesar 0.978 menunjukkan kemampuan model yang sangat baik dalam membedakan antara kelas positif dan negatif.

Semakin mendekati kurva ROC ke sudut kiri atas, semakin baik performa model dalam memisahkan

Tabel 3. Hasil Evaluasi Eksperimen Pertama TF-IDF

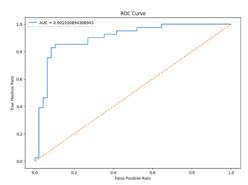
Metrik	Nilai
Akurasi	0.8764
Presisi	0.8947
Recall	0.8293
F1-score	0.8608
AUC-ROC	0.9019

Tabel 4. Hasil Evaluasi Eksperimen Kedua TF-IDF

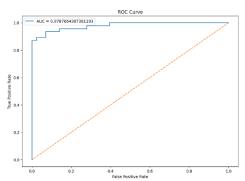
Metrik	Nilai
Akurasi	0.9326
Presisi	0.9348
Recall	0.9348
F1-score	0.9348
AUC-ROC	0.9788

Tabel 5. Hasil Evaluasi Eksperimen Ketiga TF-IDF

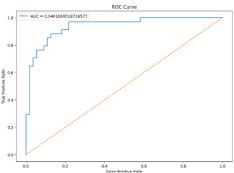
 OT OT TIME IT E THE COURT EILE	permient reenga ir ibi
Metrik	Nilai
Akurasi	0.8652
Presisi	0.7895
Recall	0.8824
F1-score	0.8333
AUC-ROC	0.9401



Gambar 2. Kurva ROC Eksperimen Pertama TF-IDF



Gambar 3. Kurva ROC Eksperimen Kedua TF-IDF



Gambar 4. Kurva ROC Eksperimen Ketiga TF-IDF

kelas positif dari kelas negatif. Pada eksperimen ketiga TF-IDF, model dievaluasi menggunakan beberapa metrik performa yaitu akurasi, presisi, recall, F1-score, dan AUC-ROC yang ditunjukkan pada Tabel 5. Akurasi model sebesar 0.865 menunjukkan bahwa model mampu memprediksi dengan benar sekitar 86.5% dari seluruh data uji yang digunakan. Presisi untuk kelas 0 adalah 0.92 dan untuk kelas 1 adalah 0.79, mengindikasikan bahwa dari seluruh prediksi positif yang dibuat oleh model, 92% untuk kelas 0 dan 79% untuk kelas 1 adalah benar. Recall untuk kelas 0 adalah 0.85 dan untuk kelas 1 adalah 0.88, yang menunjukkan bahwa model mampu menemukan 85% dari semua kasus positif sebenarnya untuk kelas 0 dan 88% untuk kelas 1. F1score, yang merupakan rata-rata harmonis antara presisi dan recall, tercatat sebesar 0.89 untuk kelas 0 dan 0.83 untuk kelas 1. Nilai AUC-ROC model sebesar 0.940 menunjukkan kemampuan model yang sangat baik dalam membedakan antara kelas positif dan negatif. Semakin mendekati kurva ROC ke sudut kiri atas, semakin baik performa model dalam memisahkan kelas positif dari kelas negatif.

3.2. Hasil Eksperimen TF-RF

Pada eksperimen pertama TF-RF, model dievaluasi menggunakan beberapa metrik performa yaitu akurasi, presisi, recall, F1-score, dan AUC-ROC yang ditunjukkan pada Tabel 6. Akurasi model sebesar 0.787 menunjukkan bahwa model mampu memprediksi dengan benar sekitar 78.7% dari seluruh data uji yang digunakan. Presisi untuk kelas nonsensitif adalah 0.76 dan untuk kelas sensitif adalah 0.82, yang mengindikasikan bahwa dari seluruh prediksi positif yang dibuat oleh model, 76% untuk kelas non-sensitif dan 82% untuk kelas sensitif adalah benar. Recall untuk kelas non-sensitif adalah 0.88 dan untuk kelas sensitif adalah 0.68, yang menunjukkan bahwa model mampu menemukan 88% dari semua kasus positif sebenarnya untuk kelas non-sensitif dan 68% untuk kelas sensitif. F1-score, yang merupakan rata-rata harmonis antara presisi dan recall, tercatat sebesar 0.82 untuk kelas non-sensitif dan 0.75 untuk kelas sensitif. Nilai AUC-ROC model sebesar 0.850 menunjukkan kemampuan model yang baik dalam membedakan antara kelas positif dan negatif.

Semakin mendekati kurva ROC ke sudut kiri atas, semakin baik performa model dalam memisahkan kelas positif dari kelas negatif. Pada eksperimen kedua TF-RF, model dievaluasi menggunakan beberapa metrik performa yaitu akurasi, presisi, recall, F1-score, dan AUC-ROC yang ditunjukkan pada Tabel 7. Akurasi model sebesar 0.854 menunjukkan bahwa model mampu memprediksi dengan benar sekitar 85.4% dari seluruh data uji yang digunakan. Presisi untuk kelas non-sensitif adalah 0.81 dan untuk kelas sensitif adalah 0.90, yang mengindikasikan bahwa dari seluruh prediksi positif yang dibuat oleh model, 81% untuk kelas non-sensitif dan 90% untuk kelas sensitif adalah benar. Recall untuk kelas non-sensitif adalah 0.91 dan untuk kelas sensitif adalah 0.80, yang menunjukkan bahwa model mampu menemukan 91% dari semua kasus positif sebenarnya untuk kelas non-sensitif dan 80% untuk kelas sensitif. F1-score, yang merupakan rata-rata harmonis antara presisi dan recall, tercatat sebesar 0.86 untuk kelas non-sensitif dan 0.85 untuk kelas sensitif. Nilai AUC-ROC model sebesar 0.903 menunjukkan kemampuan model yang sangat baik dalam membedakan antara kelas positif dan negatif.

Akurasi model sebesar 0.854 menunjukkan bahwa model mampu memprediksi dengan benar sekitar 85.4% dari seluruh data uji yang digunakan. Presisi untuk kelas non-sensitif adalah 0.89 dan untuk kelas sensitif adalah 0.80, yang mengindikasikan bahwa dari seluruh prediksi positif yang dibuat oleh model, 89% untuk kelas non-sensitif dan 80% untuk kelas sensitif adalah benar. Recall untuk kelas nonsensitif adalah 0.87 dan untuk kelas sensitif adalah 0.82, yang menunjukkan bahwa model mampu menemukan 87% dari semua kasus positif sebenarnya untuk kelas non-sensitif dan 82% untuk kelas sensitif. F1-score, yang merupakan rata-rata harmonis antara presisi dan recall, tercatat sebesar 0.88 untuk kelas non-sensitif dan 0.81 untuk kelas sensitif. Nilai AUC-ROC model sebesar 0.881 menunjukkan kemampuan model yang sangat baik dalam membedakan antara kelas positif dan negatif. mendekati kurva ROC ke sudut kiri atas, semakin baik performa model dalam memisahkan kelas positif dari kelas negatif.

3.3. Perbandingan Hasil TF-IDF dan TF-RF

melakukan Setelah tiga eksperimen menggunakan metode TF-IDF dan TF-RF, diperoleh hasil evaluasi performa model yang bervariasi. Tabel 8 merangkum hasil akurasi, presisi, recall, F1-score, dan AUC-ROC dari masing-masing metode. Dari Tabel 9, terlihat bahwa metode TF-IDF cenderung memberikan hasil yang lebih baik dibandingkan dengan metode TF-RF. Eksperimen kedua dari TF-IDF menunjukkan hasil yang paling optimal dengan nilai akurasi, presisi, recall, F1-score, dan AUC-ROC yang lebih tinggi dibandingkan dengan eksperimen lainnya. Pada metode TF-RF, eksperimen kedua juga memberikan hasil terbaik, namun tetap di bawah performa TF-IDF.

Untuk mendapatkan gambaran umum tentang performa model secara keseluruhan dari masingmasing metode, dihitung nilai rata-rata dari masingmasing metrik pada Tabel 10. Pada Tabel 10 terlihat bahwa metode TF-IDF memiliki nilai rata-rata yang lebih tinggi dibandingkan dengan TF-RF untuk semua metrik performa. Nilai rata-rata akurasi TF-IDF sebesar 0.891 menunjukkan bahwa model TF-IDF mampu memprediksi dengan benar sekitar 89.1% dari seluruh data uji, sedangkan model TF-RF hanya sebesar 83.2%. Presisi, recall, dan F1-score yang lebih tinggi pada TF-IDF juga menunjukkan bahwa metode ini lebih efektif dalam

Tabel 6. Hasil Evaluasi Eksperimen Pertama TF-RF

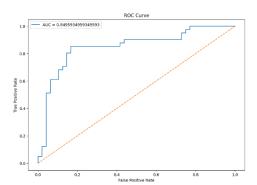
Metrik	Nilai
Akurasi	0.7865
Presisi	0.8235
Recall	0.6829
F1-score	0.7467
AUC-ROC	0.8496

Tabel 7 Hasil Evaluasi Eksperimen Kedua TE-RE

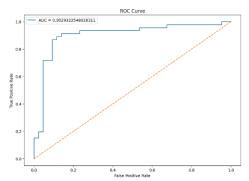
Metrik Nilai		
Nilai		
0.8539		
0.9024		
0.8043		
0.8506		
0.9029		

Tabel 8 Hasil Evaluasi Eksperimen Kedua TF-RF

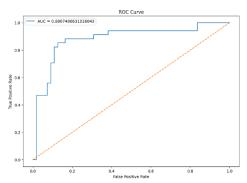
Metrik	Nilai
Akurasi	0.8539
Presisi	0.8
Recall	0.8235
F1-score	0.8116
AUC-ROC	0.8807



Gambar 5. Kurva ROC Eksperimen Pertama TF-RF



Gambar 6. Kurva ROC Eksperimen Kedua TF-RF



Gambar 7. Kurva ROC Eksperimen Ketiga TF-RF

mengklasifikasikan data dengan lebih sedikit kesalahan.

Berdasarkan hasil dari tiga eksperimen TF-IDF yang dilakukan, metode TF-IDF menunjukkan kinerja yang baik dalam mengklasifikasikan data. Eksperimen kedua memberikan hasil yang paling optimal di antara ketiga eksperimen, dengan nilai akurasi, presisi, *recall*, *F1-score*, dan *AUC-ROC* yang lebih tinggi. Nilai rata-rata dari ketiga eksperimen menunjukkan bahwa model memiliki kemampuan yang seimbang dan baik dalam mengklasifikasikan data. Secara keseluruhan, metode TF-IDF berhasil menghasilkan model dengan performa yang cukup baik, dengan nilai akurasi rata-rata sebesar 89.1% dan *AUC-ROC* rata-rata sebesar 94.0%.

Sedangkan, berdasarkan hasil dari tiga eksperimen TF-RF yang dilakukan, metode TF-RF menunjukkan kinerja yang cukup baik dalam mengklasifikasikan data. Eksperimen memberikan hasil yang paling optimal di antara ketiga eksperimen, dengan nilai akurasi, presisi, recall, F1-score, dan AUC-ROC yang lebih tinggi. Nilai rata-rata dari ketiga eksperimen menunjukkan bahwa model memiliki kemampuan yang baik dalam mengklasifikasikan data, meskipun terdapat ruang untuk perbaikan terutama pada nilai recall. Secara keseluruhan, metode TF-RF berhasil menghasilkan model dengan performa yang baik, dengan nilai akurasi rata-rata sebesar 83.2% dan AUC-ROC ratarata sebesar 87.8%.

Dari perbandingan hasil eksperimen TF-IDF dan TF-RF, dapat disimpulkan bahwa metode TF-IDF memiliki performa yang lebih baik dibandingkan metode TF-RF dalam konteks tugas klasifikasi teks. Hal ini ditunjukkan oleh nilai rata-rata akurasi, presisi, recall, F1-score, dan AUC-ROC yang lebih tinggi pada eksperimen TF-IDF dibandingkan dengan TF-RF. Keunggulan metode TF-IDF dalam mendeteksi informasi sensitif memberikan nilai strategis bagi industri jasa keuangan, terutama dalam menghadapi potensi kebocoran informasi yang dapat terjadi akibat tindakan karyawan yang akan resign. Dalam konteks ini, akurasi dan keandalan model berbasis TF-IDF memungkinkan identifikasi

informasi sensitif dengan tingkat ketepatan yang lebih tinggi, sehingga perusahaan dapat mengambil langkah preventif lebih dini untuk melindungi data penting. Dengan demikian, model berbasis TF-IDF tidak hanya unggul dalam performa teknis tetapi juga relevan secara praktis untuk meningkatkan keamanan informasi perusahaan dalam skenario dunia nyata.

4. KESIMPULAN

Penelitian ini berfokus pada pengembangan dan evaluasi model pembelajaran mesin berbasis *Convolutional Neural Networks* (CNN) untuk mendeteksi informasi sensitif dalam dokumen teks pada industri jasa keuangan, dengan membandingkan

Tabel 9. Hasil Evaluasi Eksperimen TF-IDF dan TF-RF

Metrik	,	FF-IDF			TF-RF	
			Ekspe	erimen		
	#1	#2	#3	#1	#2	#3
Akurasi	0.87	0.93	0.86	0.78	0.85	0.85
Presisi	0.89	0.93	0.78	0.82	0.90	0.80
Recall	0.82	0.93	0.88	0.68	0.80	0.82
F1-score	0.86	0.93	0.83	0.74	0.85	0.81
AUC-ROC	0.90	0.97	0.94	0.85	0.90	0.88

Tabel 10. Rata-Rata Evaluasi Eksperimen TF-IDF Dan TF-RF

Metrik	Rata-rata TF-IDF	Rata-rata TF-RF
Akurasi	0.891	0.832
Presisi	0.873	0.842
Recall	0.882	0.770
F1-score	0.876	0.803
AUC-ROC	0.940	0.878

dua metode pembobotan, yaitu *Term Frequency-Inverse Document Frequency* (TF-IDF) dan *Term Frequency-Relevance Frequency* (TF-RF). Berdasarkan hasil penelitian yang dilakukan, terdapat beberapa kesimpulan utama yang dapat diambil.

Model CNN yang diterapkan dengan metode pembobotan TF-IDF menunjukkan kinerja yang lebih baik dibandingkan dengan TF-RF dalam hal akurasi, presisi, recall, dan F1-score. Hasil eksperimen menunjukkan bahwa pembobotan TF-IDF lebih efektif dalam meningkatkan kinerja model CNN dalam mendeteksi informasi sensitif. Hasil ini mengindikasikan bahwa metode TF-IDF mampu memberikan bobot yang lebih tepat pada kata-kata yang signifikan dalam dokumen, sehingga model dapat mengenali pola-pola yang relevan dengan lebih akurat. Model CNN yang sudah dibangun memiliki keunggulan sepert CNN efektif dalam menangkap pola dalam data teks, misalnya frasa atau kombinasi kata yang relevan untuk klasifikasi sensitive dan non sensitif. Kemudian Dengan TF-IDF, model dapat lebih fokus pada kata-kata yang relevan untuk mendeteksi informasi sensitif, terutama dalam dokumen teks yang panjang. Hal yang menjadi perhatian adalah penggunakan dataset yang memadai agar model tidak mengalami overfitting.

Selain itu, penelitian ini juga menegaskan pentingnya pemilihan metode pembobotan yang tepat dalam pengembangan model pembelajaran mesin. Metode TF-IDF terbukti lebih unggul dalam konteks deteksi informasi sensitif pada dokumen teks, memberikan keseimbangan yang baik antara presisi dan recall serta mengurangi kesalahan prediksi positif. Hasil ini memberikan dasar yang kuat untuk implementasi model CNN dengan pembobotan TF-IDF dalam industri jasa keuangan guna meningkatkan keamanan informasi dan mencegah potensi penyalahgunaan data oleh karyawan yang akan mengundurkan diri.

DAFTAR PUSTAKA

- AHSAN, M.M., MAHMUD, M.A.P., SAHA, P.K., GUPTA, K.D. dan SIDDIQUE, Z., 2021. Effect of Data Scaling Methods on Machine Learning Algorithms and Model Performance. Technologies, 9(3), p.52.
- AKBAR, A.T., SAIFULLAH, S., PRAPCOYO, H., PEMBANGUNAN, U., VETERAN, N., SLEMAN, K. dan KORESPONDENSI, P., 2024. Klasifikasi Ekspresi Wajah Menggunakan Covolutional Neural Network. Jurnal Teknologi Informasi dan Ilmu Komputer, 11(6), pp.1399-1412.
- AL-SHEHARI, T. dan ALSOWAIL, R.A., 2021, An insider data leakage detection using one-hot encoding, synthetic minority oversampling and machine learning techniques. Entropy, 23(10).
- ALI, A., ABD RAZAK, S., OTHMAN, S.H., EISA, T.A.E., AL-DHAQM, A., NASSER, M., ELHASSAN, T., ELSHAFIE, H. dan SAIF, A., 2022. Financial Fraud Detection Based on Machine Learning: A Systematic Review. Applied Sciences Literature (Switzerland), 12(19).
- ALI, B.H., JALAL, A.A. dan IBRAHEM AL-OBAYDY, W.N., 2020. Data loss prevention by using MRSH-v2 algorithm. International Journal of Electrical and Computer Engineering, 10(4), pp.3615-3622.
- ALZUBAIDI, L., ZHANG, J., HUMAIDI, A.J., AL-DUJAILI, A., DUAN, Y., AL-SHAMMA, O., SANTAMARÍA, J., FADHEL, M.A., AL-AMIDIE, M. dan FARHAN, L., 2021. Review of deep learning: concepts, CNN challenges, applications, architectures, future directions. Journal of Big Data, p.53. Available [online] 8. https://doi.org/10.1186/s40537-021- 00444-8> [Accessed 4 Jun. 2021].
- INDRADEWI, DANANJAYA, I.K.W. dan I.G.A.A.D., 2023. Perbandingan Metode Pembobotan TF-RF Dan TF-ABS Pada Kategorisasi Berita Di BDI Denpasar. *SINTECH* (Science and Information Technology) Journal, 6(1), pp.16–25.
- GHORBANI, R. dan GHOUSI, R., 2020. Comparing

- Different Resampling Methods in Predicting Students' Performance Using Machine Learning Techniques. IEEE Access, 8, pp.67899–67911.
- GUHA, A., SAMANTA, D., BANERJEE, A. dan AGARWAL, D., 2021. A Deep Learning Model for Information Loss Prevention from Multi-Page Digital Documents. IEEE Access, 9, pp.80451–80465.
- HARMANDINI, K.P. dan L, K.M., 2024. Analysis of TF-IDF and TF-RF Feature Extraction on Product Review Sentiment. Sinkron, 8(2), pp.929-937.
- IRENA, B. dan SETIAWAN, E.B., 2020. Fake News (Hoax) Identification on Social Media Twitter using Decision Tree C4.5 Method. Jurnal RESTI (Rekayasa Sistem dan Teknologi Informasi), 4(4), pp.711–716.
- KRISNABAYU, R.Y., RIDOK, A. dan BUDI, A.S., 2021. Hepatitis Detection using Random Forest based on SVM-RFE (Recursive Feature Elimination) Feature Selection and SMOTE. 6th International Conferenceon Sustainable Information Engineering and Technology 2021 (SIET '21), pp.151-156.
- MAHARDIKA, M.A.R., YUDISTIRA, N. dan RIDOK, A., 2024. Sistem Rekognisi Citra Digital Bahasa Isyarat Menggunakan Convolutional Neural Network dan Spatial Transformer. Jurnal Teknologi Informasi dan Ilmu Komputer, 11(6), pp.11590-1168.
- NUGROHO, K.S., ISTIADI, I. dan MARISA, F., 2020. Naive Bayes classifier optimization for text classification on e-government using particle swarm optimization. Jurnal Teknologi dan Sistem Komputer, 8(1), pp.21–26.
- NURDIAWAN, O., SUSANA, H. dan FAQIH, A., 2024. Deep Learning for Polycystic Ovarian Syndrome Classification Using Neural Network. JITK Convolutional (Jurnal Ilmu Pengetahuan dan Teknologi Komputer), 9(2), pp.218–226.
- RACHMAN, F.P. dan SANTOSO, H., 2021. Perbandingan Model Deep Learning untuk Klasifikasi Sentiment Analysis dengan Teknik Natural Languange Processing. Jurnal Teknologi dan Manajemen *Informatika*, 7(2), pp.113–121.
- RAMADHAN, R., SARI, Y.A. dan ADIKARA, P.P., 2021. Perbandingan Pembobotan Term Frequency-Inverse Document Frequency dan Term Frequency-Relevance Frequency terhadap Fitur N-Gram pada Analisis Sentimen. Jurnal Pengembangan Teknologi Informasi dan Ilmu Komputer, [online] pp.5075–5079. Available https://j-ptiik.ub.ac.id/index.php/j-

- ptiik/article/view/10173>.
- SARI, Y., BASKARA, A.R., PRAKOSO, P.B., ROYANI, N., MANGKURAT, U.L. dan KORESPONDENSI, P., 2022. Perbandingan Metode Pembobotan Tf-Rf Dan Tf-Idf Dengan Dikombinasikan Dengan Weighted Tree Similarity Comparison of Tf-Rf and . Weighting Methods Tf-Idf With Combination With Weighted Tree Similarity for a Book. *Jurnal Teknologi Informasi dan Ilmu Komputer*, 9(3).
- SHAKTI, F.N. dan HIDAYANTO, A.N., 2024.

 Measurement of Employee Information
 Security Awareness: Case Study At
 Financial Institution. *JITK (Jurnal Ilmu Pengetahuan dan Teknologi Komputer)*,
 9(2), pp.172–179.
- SOLDATOS, J. dan KYRIAZIS, D., 2022. Big Data and Artificial Intelligence in Digital Finance. Big Data and Artificial Intelligence in Digital Finance.
- STOJANOVIĆ, B., BOŽI, J., HOFER-SCHMITS, K., NAHRGANG, K., WEBER, A., BADII, A., SUNDARAM, M., JORDAN, E. dan RUNEVIC, J., 2021. Follow the Trail: Machine Learning for Fraud Detection in Fintech Applications. *Sensors*, [online] 21(5), p.1594. Available at: https://doi.org/s21051594>.
- SUARTANA, I.M., 2022. Analisis Penerapan Deep

- Learning untuk Klasifikasi Serangan Terhadap Keamanan Jaringan. *Klik-Kumpulan Jurnal Ilmu Komputer*, 9(1), pp.100–109.
- TAJBAKHSH, N., SHIN, J.Y., GURUDU, S.R., HURST, R.T., KENDALL, C.B., GOTWAY, M.B. dan LIANG, J., 2017. Convolutional Neural Networks for Medical Image Analysis: Full Training or Fine Tuning? *IEEE Transactions on Medical Imaging*, 35(5), pp.1299–1312.
- TAMA, F.R. dan SIBARONI, Y., 2023. Fake News (Hoaxes) Detection on Twitter Social Media Content through Convolutional Neural Network (CNN) Method. *JINAV: Journal of Information and Visualization*, 4(1), pp.70–78.
- TEMPL, M. dan SARIYAR, M., 2022. A systematic overview on methods to protect sensitive data provided for various analyses. *International Journal of Information Security*, [online] 21(6), pp.1233–1246. Available at: https://doi.org/10.1007/s10207-022-00607-5>.
- YUHANA, U.L., IMAMAH, I., FATICHAH, C. dan SANTOSO, B.J., 2022. Effectiveness of Deep Learning Approach for Text Classification in Adaptive Learning. *Jurnal Ilmiah Kursor*, 11(3), p.137.